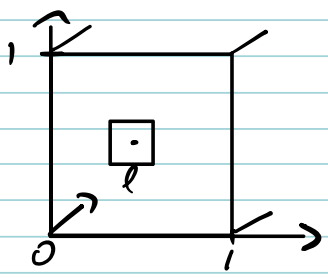


Curse of Dimensionality

Assume $x_i \in [0,1]^d$ (i.e. the d dimensional unit hypercube).
 All data is drawn uniformly at random.
 Let $k=10$.

Let l be the edge length of the smallest hypercube that contains all k nearest neighbors of a test point x .



$$l^d \approx \frac{k}{n} \Rightarrow l \approx \left(\frac{k}{n}\right)^{1/d}$$

↑ volume of mini cube containing the k neighbors
 ↑ Total volume of hypercube $[0,1]^d$ is $1^d=1$
 $\frac{k}{n}$ is the fraction that k points take up. (because points are uniformly sampled)

If $n=1000$ how big is l ?

d	2	10	100	1000
l	0.1	0.63	0.955	0.9954

Almost the entire space is needed to fit 10 nearest neighbors.

This means nearest neighbors are not similar, violating the k -NN assumption!

How many points would we need for l to be small?
 Fix $l=0.1$

$$l^d = \frac{k}{n} \Rightarrow n = k \left(\frac{1}{l}\right)^d = k 10^d \leftarrow \text{grows exponentially with } d!$$

Rescue to the curse:

Data can have structure:

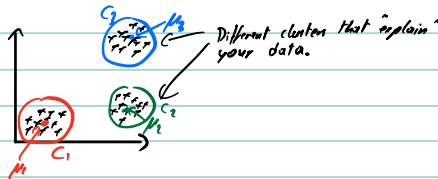
- Data can lie on intrinsically low dimensional subspaces or sub-manifolds.
- Data can be clustered (very non-uniform).

Unsupervised Learning

Given: Unlabeled data $D = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$ ← No labels y_i
 Unsupervised learning attempts to find (hidden) structure in your data.

Clustering: K-Means

Assumptions: - We know there are k clusters
 - Similar points share the same cluster.



Let data indices be $\{1, \dots, n\}$. Let $C_l \subseteq \{1, \dots, n\}$ be all points that fall into the l -th cluster.
 Every point falls exactly into one cluster: $C_1 \cup C_2 \cup \dots \cup C_k = \{1, \dots, n\}$
 (Clusters are mutually exclusive.) $C_l \cap C_{l'} = \emptyset \quad \forall l \neq l'$

Cluster center: $\vec{\mu}_l = \frac{1}{|C_l|} \sum_{i \in C_l} \vec{x}_i$

Loss function: $\mathcal{L}(C_1, \dots, C_k) = \frac{1}{n} \sum_{l=1}^k \sum_{i \in C_l} \|\vec{x}_i - \vec{\mu}_l\|^2$ ← Average distance from each point to its cluster center.

↳ If you were to represent each \vec{x}_i with its corresponding $\vec{\mu}_l$ instead, this would be the avg. squared error.

① Minimize \mathcal{L} if C_1, \dots, C_k are fixed: $\vec{\mu}_l = \underset{\vec{\mu}_l}{\operatorname{argmin}} \mathcal{L}(C_1, \dots, C_k) \Rightarrow \vec{\mu}_l = \frac{1}{|C_l|} \sum_{i \in C_l} \vec{x}_i$

(Assign μ_1, \dots, μ_k)

For each μ_l , the contribution to the loss is $\sum_{i \in C_l} (x_i - \mu_l)^2$ ← Parabola

mean of all $x_i \in C_l$

② Minimize \mathcal{L} if μ_1, \dots, μ_k are fixed: $C_l = \underset{C_l \subseteq \{1, \dots, n\}}{\operatorname{argmin}} \sum_{i \in C_l} \|\vec{x}_i - \mu_l\|^2$

(Assign C_1, \dots, C_k)

↑
All points that are closest to μ_l .

For each x_i you can contribute $(x_i - \mu_1)^2$ or $(x_i - \mu_2)^2$ or \dots or $(x_i - \mu_k)^2$.
 For each x_i pick the smallest one to minimize \mathcal{L} .

K-means algorithm: 1) Initialize $\vec{\mu}_1, \dots, \vec{\mu}_k$ somehow
 (aka Lloyd's algorithm) (e.g. random points, or heuristic)

2) Assign C_1, \dots, C_k using ②

3) Assign $\vec{\mu}_1, \dots, \vec{\mu}_k$ using ①

Repeat until convergence

Easy to show: \mathcal{L} can never increase.

Initialization: - Can be unucky. Perform multiple runs, pick solution with lowest \mathcal{L} .

- k-means++ (Arthur et al. 2009): - pick μ_1 randomly from D

- For $l=2$ to k :

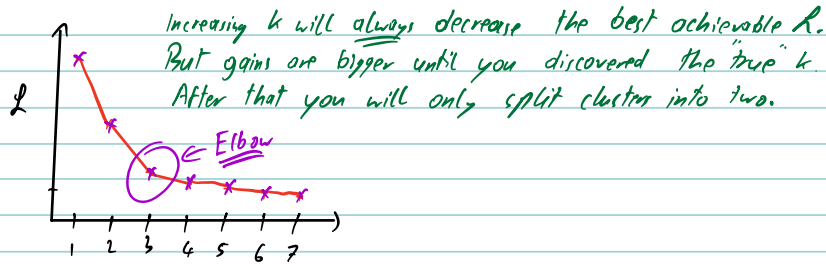
- $\forall x_i \in D$ define $d_i = \min_j (x_i - \mu_j)^2$

- pick μ_l randomly from D proportional to d_i^2 .

↳ Avoids cluster centers that are very close

How to find k if it is unknown?

Run k -means for $k=1,2,\dots$ monitor the loss L
 (For each k you have to average over multiple runs, due to the randomness of the initialization.)



Strengths:

- Easy to implement
- fast for small k
- easily parallelizable

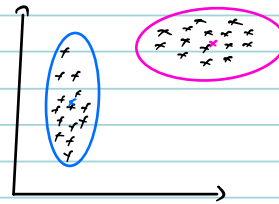
Limitations:

- sensitive to outlier
- slow for large k
- assumes all clusters are spherical
- makes only hard cluster assignments

Generalization: Gaussian Mixture Model (GMM)

Each cluster becomes a Gaussian $N(\mu_c, \Sigma_c)$

Almost the same as k -means, but each point has a distribution over clusters $\gamma_i = P(c_c | x_i)$.



$$P(c_c | \tilde{x}_i) = \frac{P(\tilde{x}_i | c_c) P(c_c)}{\sum_j P(\tilde{x}_i | c_j) P(c_j)} \quad (\text{Bayes Rule})$$

$$\vec{\mu}_c \leftarrow \frac{1}{\sum_j \gamma_{ij}} \sum_{i=1}^n \gamma_{ij} \tilde{x}_i \leftarrow \text{weighted average / expected center of cluster } c$$

$$\vec{\Sigma}_c \leftarrow \frac{1}{\sum_j \gamma_{ij}} \sum_{i=1}^n \gamma_{ij} (\tilde{x}_i - \mu_c)(\tilde{x}_i - \mu_c)^T$$

$$\pi_c \leftarrow \frac{\sum_i \gamma_{ij}}{n} \leftarrow \pi_c = P(c_c)$$

$$\gamma_{ij} \leftarrow \frac{P(x_i; \mu_c, \Sigma_c) \pi_c}{\sum_{c'} P(x_i; \mu_{c'}, \Sigma_{c'}) \pi_{c'}}$$

M-Step

Maximization

E-Step

Estimation

Strengths:

- more flexible than k -means
- probabilities tell you how well a sample fits into cluster

Weaknesses:

- slower than k -means
- can suffer from singularities (single point clusters)

(If $\Sigma_c = I\sigma^2$ GMM becomes k -means as $\sigma^2 \rightarrow 0$.)

spherical

hard assignments